# NAG Library Routine Document

# G02CDF

## 1    Purpose

G02CDF performs a simple linear regression with no constant, with dependent variable $y$ and independent variable $x$, omitting cases involving missing values.

## 2    Specification

```
SUBROUTINE G02CDF (N, X, Y, XMISS, YMISS, RESULT, IFAIL)

INTEGER          N, IFAIL
REAL (KIND=nag_wp) X(N), Y(N), XMISS, YMISS, RESULT(21)
```

## 3    Description

G02CDF fits a straight line of the form

$$y = bx$$

to those of the data points

$$(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$$

that do not include missing values, such that

$$y_i = bx_i + e_i$$

for those $(x_i, y_i)$, for $i = 1, 2, \ldots, n$    $(n \geq 2)$ which do not include missing values.

The routine eliminates all pairs of observations $(x_i, y_i)$ which contain a missing value for either $x$ or $y$, and then calculates the regression coefficient, $b$, and various other statistical quantities by minimizing the sum of the $e_i^2$ over those cases remaining in the calculations.

The input data consists of the $n$ pairs of observations $(x_1, y_1), (x_2, y_2), \ldots, (x_n, y_n)$ on the independent variable $x$ and the dependent variable $y$.

In addition two values, *xm* and *ym*, are given which are considered to represent missing observations for $x$ and $y$ respectively. (See Section 7).

Let $w_i = 0$, if the $i$th observation of either $x$ or $y$ is missing, i.e., if $x_i = xm$ and/or $y_i = ym$; and $w_i = 1$ otherwise, for $i = 1, 2, \ldots, n$.

The quantities calculated are:

(a)    Means:

$$\bar{x} = \frac{\sum\limits_{i=1}^{n} w_i x_i}{\sum\limits_{i=1}^{n} w_i}; \qquad \bar{y} = \frac{\sum\limits_{i=1}^{n} w_i y_i}{\sum\limits_{i=1}^{n} w_i}.$$

(b)    Standard deviations:

$$s_x = \sqrt{\frac{\sum\limits_{i=1}^{n} w_i(x_i - \bar{x})^2}{\sum\limits_{i=1}^{n} w_i - 1}}; \qquad s_y = \sqrt{\frac{\sum\limits_{i=1}^{n} w_i(y_i - \bar{y})^2}{\sum\limits_{i=1}^{n} w_i - 1}}.$$

(c)    Pearson product-moment correlation coefficient:

$$r = \frac{\sum\limits_{i=1}^{n} w_i(x_i - \bar{x})(y_i - \bar{y})}{\sum\limits_{i=1}^{n} w_i(x_i - \bar{x})^2 \sum\limits_{i=1}^{n} w_i(y_i - \bar{y})^2}.$$

(d)    The regression coefficient, $b$:

$$b = \frac{\sum\limits_{i=1}^{n} w_i x_i y_i}{\sum\limits_{i=1}^{n} w_i x_i^2}.$$

(e)    The sum of squares attributable to the regression, $SSR$, the sum of squares of deviations about the regression, $SSD$, and the total sum of squares, $SST$:

$$SST = \sum\limits_{i=1}^{n} w_i y_i^2; \qquad SSD = \sum\limits_{i=1}^{n} w_i(y_i - bx_i)^2; \qquad SSR = SST - SSD.$$

(f)    The degrees of freedom attributable to the regression, $DFR$, the degrees of freedom of deviations about the regression, $DFD$, and the total degrees of freedom, $DFT$:

$$DFT = \sum\limits_{i=1}^{n} w_i; \qquad DFD = \sum\limits_{i=1}^{n} w_i - 1; \qquad DFR = 1.$$

(g)    The mean square attributable to the regression, $MSR$, and the mean square of deviations about the regression, $MSD$:

$$MSR = SSR/DFR; \qquad MSD = SSD/DFD.$$

(h)    The $F$ value for the analysis of variance:

$$F = MSR/MSD.$$

(i)    The standard error of the regression coefficient:

$$se(b) = \sqrt{\frac{MSD}{\sum\limits_{i=1}^{n} w_i x_i^2}}.$$

(j)    The $t$ value for the regression coefficient:

$$t(b) = \frac{b}{se(b)}.$$

(k)    The number of observations used in the calculations:

$$n_c = \sum\limits_{i=1}^{n} w_i.$$

# 4    References

Draper N R and Smith H (1985) *Applied Regression Analysis* (2nd Edition) Wiley

## 5 Parameters

1:     N – INTEGER                *Input*

*On entry*: $n$, the number of pairs of observations.

*Constraint*: $N \geq 2$.

2:     X(N) – REAL (KIND=nag_wp) array               *Input*

*On entry*: $X(i)$ must contain $x_i$, for $i = 1, 2, \ldots, n$.

3:     Y(N) – REAL (KIND=nag_wp) array               *Input*

*On entry*: $Y(i)$ must contain $y_i$, for $i = 1, 2, \ldots, n$.

4:     XMISS – REAL (KIND=nag_wp)               *Input*

*On entry*: the value $xm$, which is to be taken as the missing value for the variable $x$ (see Section 7).

5:     YMISS – REAL (KIND=nag_wp)               *Input*

*On entry*: the value $ym$, which is to be taken as the missing value for the variable $y$ (see Section 7).

6:     RESULT(21) – REAL (KIND=nag_wp) array              *Output*

*On exit*: the following information:

| | |
|---|---|
| RESULT(1) | $\bar{x}$, the mean value of the independent variable, $x$; |
| RESULT(2) | $\bar{y}$, the mean value of the dependent variable, $y$; |
| RESULT(3) | $s_x$, the standard deviation of the independent variable, $x$; |
| RESULT(4) | $s_y$, the standard deviation of the dependent variable, $y$; |
| RESULT(5) | $r$, the Pearson product-moment correlation between the independent variable $x$ and the dependent variable, $y$; |
| RESULT(6) | $b$, the regression coefficient; |
| RESULT(7) | the value 0.0; |
| RESULT(8) | $se(b)$, the standard error of the regression coefficient; |
| RESULT(9) | the value 0.0; |
| RESULT(10) | $t(b)$, the $t$ value for the regression coefficient; |
| RESULT(11) | the value 0.0; |
| RESULT(12) | $SSR$, the sum of squares attributable to the regression; |
| RESULT(13) | $DFR$, the degrees of freedom attributable to the regression; |
| RESULT(14) | $MSR$, the mean square attributable to the regression; |
| RESULT(15) | $F$, the $F$ value for the analysis of variance; |
| RESULT(16) | $SSD$, the sum of squares of deviations about the regression; |
| RESULT(17) | $DFD$, the degrees of freedom of deviations about the regression; |
| RESULT(18) | $MSD$, the mean square of deviations about the regression; |
| RESULT(19) | $SST$, the total sum of squares |
| RESULT(20) | $DFT$, the total degrees of freedom; |
| RESULT(21) | $n_c$, the number of observations used in the calculations. |

7:     IFAIL – INTEGER               *Input/Output*

*On entry*: IFAIL must be set to 0, $-1$ or 1. If you are unfamiliar with this parameter you should refer to Section 3.3 in the Essential Introduction for details.

For environments where it might be inappropriate to halt program execution when an error is detected, the value $-1$ or 1 is recommended. If the output of error messages is undesirable, then the value 1 is recommended. Otherwise, if you are not familiar with this parameter, the recommended value is 0. **When the value $-1$ or 1 is used it is essential to test the value of IFAIL on exit.**

*On exit*: IFAIL $= 0$ unless the routine detects an error or a warning has been flagged (see Section 6).

## 6 Error Indicators and Warnings

If on entry IFAIL $= 0$ or $-1$, explanatory error messages are output on the current error message unit (as defined by X04AAF).

Errors or warnings detected by the routine:

IFAIL $= 1$

On entry, $N < 2$.

IFAIL $= 2$

After observations with missing values were omitted, fewer than two cases remained.

IFAIL $= 3$

After observations with missing values were omitted, all remaining values of at least one of the variables $x$ and $y$ were identical.

## 7 Accuracy

G02CDF does not use *additional precision* arithmetic for the accumulation of scalar products, so there may be a loss of significant figures for large $n$.

You are warned of the need to exercise extreme care in your selection of missing values. G02CDF treats all values in the inclusive range $\left(1 \pm 0.1^{(\text{X02BEF}-2)}\right) \times xm_j$, where $xm_j$ is the missing value for variable $j$ specified in XMISS.

You must therefore ensure that the missing value chosen for each variable is sufficiently different from all valid values for that variable so that none of the valid values fall within the range indicated above.

If, in calculating $F$ or $t(b)$ (see Section 3), the numbers involved are such that the result would be outside the range of numbers which can be stored by the machine, then the answer is set to the largest quantity which can be stored as a real variable, by means of a call to X02ALF.

## 8 Further Comments

The time taken by G02CDF depends on $n$ and the number of missing observations.

The routine uses a two-pass algorithm.

## 9 Example

This example reads in eight observations on each of two variables, and then performs a simple linear regression with no constant, with the first variable as the independent variable, and the second variable as the dependent variable, omitting cases involving missing values (0.0 for the first variable, 99.0 for the second). Finally the results are printed.

### 9.1 Program Text

```
    Program g02cdfe

!     G02CDF Example Program Text

!     Mark 24 Release. NAG Copyright 2012.

!     .. Use Statements ..
      Use nag_library, Only: g02cdf, nag_wp
!     .. Implicit None Statement ..
      Implicit None
!     .. Parameters ..
      Integer, Parameter                :: nin = 5, nout = 6
!     .. Local Scalars ..
```

```
      Real (Kind=nag_wp)                   :: xmiss, ymiss
      Integer                              :: i, ifail, n
!     .. Local Arrays ..
      Real (Kind=nag_wp)                   :: reslt(21)
      Real (Kind=nag_wp), Allocatable :: x(:), y(:)
!     .. Executable Statements ..
      Write (nout,*) 'G02CDF Example Program Results'
      Write (nout,*)

!     Skip heading in data file
      Read (nin,*)

!     Read in problem size
      Read (nin,*) n

      Allocate (x(n),y(n))

!     Read in data
      Read (nin,*)(x(i),y(i),i=1,n)

!     Read in missing value flags
      Read (nin,*) xmiss, ymiss

!     Display data
      Write (nout,*) ' Case     Independent     Dependent'
      Write (nout,*) 'number     variable       variable'
      Write (nout,*)
      Write (nout,99999)(i,x(i),y(i),i=1,n)
      Write (nout,*)

!     Fit linear regression model
      ifail = 0
      Call g02cdf(n,x,y,xmiss,ymiss,reslt,ifail)

!     Display results
      Write (nout,99998) 'Mean of independent variable           = ', &
        reslt(1)
      Write (nout,99998) 'Mean of   dependent variable           = ', &
        reslt(2)
      Write (nout,99998) 'Standard deviation of independent variable = ', &
        reslt(3)
      Write (nout,99998) 'Standard deviation of   dependent variable = ', &
        reslt(4)
      Write (nout,99998) 'Correlation coefficient                = ', &
        reslt(5)
      Write (nout,*)
      Write (nout,99998) 'Regression coefficient                 = ', &
        reslt(6)
      Write (nout,99998) 'Standard error of coefficient          = ', &
        reslt(8)
      Write (nout,99998) 't-value for coefficient                = ', &
        reslt(10)
      Write (nout,*)
      Write (nout,*) 'Analysis of regression table :-'
      Write (nout,*)
      Write (nout,*) &
        '      Source        Sum of squares D.F.    Mean square     F-value'
      Write (nout,*)
      Write (nout,99997) 'Due to regression', reslt(12:15)
      Write (nout,99997) 'About  regression', reslt(16:18)
      Write (nout,99997) 'Total            ', reslt(19:20)
      Write (nout,*)
      Write (nout,99996) 'Number of cases used = ', reslt(21)

99999 Format (1X,I4,2F15.4)
99998 Format (1X,A,F8.4)
99997 Format (1X,A,F14.4,F8.0,2F14.4)
99996 Format (1X,A,F3.0)
    End Program g02cdfe
```

## 9.2 Program Data

```
G02CDF Example Program Data
8          :: N
1.0       20.0
0.0       15.5
4.0       28.3
7.5       45.0
2.5       24.5
0.0       10.0
10.0      99.0
5.0       31.2  :: End of X, Y
0.0       99.0  :: XMISS, YMISS
```

## 9.3 Program Results

```
G02CDF Example Program Results

 Case      Independent      Dependent
number      variable        variable

   1          1.0000         20.0000
   2          0.0000         15.5000
   3          4.0000         28.3000
   4          7.5000         45.0000
   5          2.5000         24.5000
   6          0.0000         10.0000
   7         10.0000         99.0000
   8          5.0000         31.2000

Mean of independent variable          =    4.0000
Mean of   dependent variable          =   29.8000
Standard deviation of independent variable =   2.4749
Standard deviation of   dependent variable =   9.4787
Correlation coefficient               =    0.9799

Regression coefficient                =    6.5833
Standard error of coefficient         =    0.8046
t-value for coefficient               =    8.1816

Analysis of regression table :-

      Source       Sum of squares  D.F.   Mean square      F-value

Due to regression    4528.9493     1.     4528.9493        66.9392
About   regression    270.6307     4.       67.6577
Total                4799.5800     5.

Number of cases used =  5.
```

_____